

Semantic Parsing for Conversational Question Answering over Knowledge Graphs

Laura Perez-Beltrachini¹, Parag Jain¹, Emilio Monti², Mirella Lapata¹

School of Informatics, University of Edinburgh¹

Amazon United Kingdom²

{lperez,parag.jain,mlap}@ed.ac.uk

monti@amazon.co.uk

Abstract

In this paper we study conversational question answering (QA) over general purpose knowledge graphs (KGs) with very large vocabularies covering thousands of concept names and relations, and millions of entities. We are interested in developing semantic parsers which understand natural language questions embedded in a conversation with a user and ground them to formal queries over definitions in a KG. To this end, we develop a dataset where user questions are annotated with SPARQL parses and system answers correspond to execution results thereof. We implement two different semantic parsing approaches and highlight the challenges of the task: dealing with large vocabularies, modelling conversation context, predicting queries with multiple entities, and generalising to new user questions at test time. We hope our dataset will serve as useful testbed for the development of conversational semantic parsers.¹

1 Introduction

Conversational information seeking is the process of acquiring information through conversations (Zamani et al., 2022). Recent years have seen an increasing number of applications aiming to satisfy users’ needs through conversational interfaces based on information retrieval (Radlinski and Craswell, 2017) and user recommendation (Janach et al., 2021). The popularity of intelligent voice assistants such as Amazon’s Alexa or Apple’s Siri has further stimulated research on question answering over general purpose knowledge graphs (e.g., Wikidata, Google’s knowledge graph). Key to question answering in this context is the ability to ground natural language onto concepts, entities, and relations in order to produce an *executable* query (e.g., SPARQL) which will retrieve an answer or *denotation* from the knowledge graph (KG).

This grounding process, known as *semantic parsing* has been either studied in the context of one or few domain-specific databases (Yu et al., 2019a; Jain and Lapata, 2021; Suhr et al., 2018) or without taking the conversational nature of the task into account (Reddy et al., 2014; Yih et al., 2016; Dubey et al., 2019; Gu et al., 2021). Conversational semantic parsing over large KGs requires handling very large vocabularies covering thousands of concept names and relations, and millions of entities rather than specialized terms consisting of hundreds of tables and column names. Moreover, information seeking conversations are by nature incremental involving interrelated rather than isolated questions. To date there is no conversational semantic parsing dataset over large KGs. Collecting natural information seeking conversations where user questions are annotated with executable queries against a KG is an expensive endeavor.

In this work, we create SPICE, a Semantic ParsIng dataset for Conversational quEstion answering over Wikidata. Our dataset consists of user-assistant interactions where natural language questions in user turns are paired with SPARQL parses and answers provided by the system correspond to SPARQL execution results. We derive this dataset from CSQA (Saha et al., 2018), an existing benchmark originally proposed for retrieval-based conversational question answering (Lan et al., 2021). CSQA consists of information seeking conversations as a sequence of user natural language questions and assistant answers (i.e., no executable queries). Conveniently, CSQA includes conversational phenomena such as coreference, ellipsis, and topic change as well as different types of questions exemplifying varying intents. Table 1 shows a conversation from SPICE and illustrates how questions (utterances on the left) are annotated with SPARQL queries (SP on the right blue box). We manually specify SPARQL templates for CSQA user questions intents, i.e., queries with entity, re-

¹Our dataset and models are released at <https://github.com/EdinburghNLP/SPICE>.

lation, and class symbols under-specified which we subsequently fill automatically to generate full SPARQL queries. We create a large-scale dataset with 197k conversations.

User questions in CSQA have been previously associated with logical forms valid under custom designed grammars (Guo et al., 2018; Kacupaj et al., 2021; Marion et al., 2021). Each semantic parser thus operates on different sets of grammar rules which hinders the comparison of proposed approaches. Different grammars may have different coverage and the semantics of different grammar terminal symbols may encapsulate different degrees of execution complexity. To promote conversational semantic parsing comparability we choose the standard SPARQL grammar.² This choice also provides generality as our proposed dataset could be further extended with new question intents without the need to redefine a custom grammar and its execution engine.

Generalisation and adaptation to new entities and concepts is key in conversational information seeking interfaces. We create different data splits where new question intents appear only at test time (Finegan-Dollak et al., 2018). With regard to modeling, we establish two strong baselines for our semantic parsing derived from existing approaches. These tackle the large vocabulary problem and the prediction of logical forms in different ways. The first approach uses *dynamic vocabularies* derived from KG subgraphs for each question and a simple sequence-to-sequence architecture to predict *complete* SPARQL queries. The other approach predicts SPARQL *query templates* and then fills in entity, relation, and type slots by means of an entity and ontology classifier. Experiments unveil different shortcomings of these approaches which range from encoding large sets of KG elements to the inability to output the same entity several times. Detailed analysis reveals that both approaches are unable to resolve coreference when the referent appears in the conversation context beyond the previous turn, struggle to resolve ellipsis, and have reduced performance on questions with multiple entities. They also have difficulty with unseen question intents, even though different operators have been seen during training. We discuss these challenges in detail and outline research directions for improving conversational semantic parsing.

2 The SPICE Dataset

The CSQA dataset (Saha et al., 2018) aims to facilitate the development of QA systems that handle complex and inter-related questions over a knowledge graph (KG). In contrast to simple factual questions that can be answered with a single KG triple (i.e., {subject, relation, object}), complex questions require manipulating sets of triples and reasoning over these. For instance, in a conversation excerpt like the one in Table 1, a question like *How many sports teams participated in that tournament?* requires numerical reasoning. Answering the question in turn \mathcal{T}_2 relies on correctly interpreting \mathcal{T}_1 .

Questions and answers in this dataset were elicited from human experts playing user and system roles as well as from crowd-workers. In a second stage, templates derived from the human-authored QA pairs were used to automatically augment the dataset. Human experts also suggested complex reasoning questions and derived templates thereof. Conversations were built from the bank of QA pairs as sequences of QA pairs exploring paths in the KG. By construction, the QA pairs in a conversation are connected through one or several entities in the KG. Questions fall into two coarse categories, *simple* and *reasoning*-based, and the way QA pairs are organised in a sequence introduces various *conversational phenomena* which we summarize below.

Simple Questions are factoid questions, seeking information related to an entity (e.g., *Which tournament did Detroit Tigers participate in?* in Table 1) or set of entities (e.g., *What are the countries of those sports teams?*).

Reasoning Questions are complex questions which require the application of numerical and logical operators over sets of entities. For instance, to answer the question *How many sports teams participated in that tournament?* requires finding the set of sports teams that participated in a given tournament (e.g., *1909 World Series*) and taking its count. Questions in this category also involve General Entities (GE) such as *tournament*, in addition to Named Entities (NE), and multiple entities (both NE and GE) in a single question (e.g., *Which tournaments have less number of participating sports teams than 1909 World Series?*). Some question types also combine multiple reasoning operators.

²<https://www.w3.org/TR/sparql11-query/>

Utterances	Annotations	Actions and Semantic Parses
\mathcal{T}_1 U: Which tournament did Detroit Tigers participate in? S: 1909 World Series	INTENT =Simple Question Single Entity ENT =[Q650855 (Detroit Tigers)], REL =[P1923 (participating team)], TYP =[Q500834 (tournament)], TRIPLE =[(Q500834,P1923,Q650855)], GOLD =[Q846847 (1909 World Series)]	AS: [filter_type, find_rev, Q650855, P1923, Q500834] SP: SELECT ?x WHERE { ?x wdt:P1923 wd:Q650855. ?x wdt:P31 wd:Q500834.}
\mathcal{T}_2 U: Which sports team was the champion of that tournament? S: Pittsburgh Pirates	INTENT =Simple Question Single Entity Indirect ENT =[Q846847 (1909 World Series)], REL =[P1346 (winner)], TYP =[Q12973014 (sports team)], TRIPLE =[(Q846847,P1346,Q12973014)], GOLD =[Q7199360 (Pittsburgh Pirates)]	AS: [filter_type, find, Q846847, P1346, Q12973014] SP: SELECT ?x WHERE { wd:Q846847 wdt:P1346 ?x. ?x wdt:P31 wd:Q12973014.}
\mathcal{T}_3 U: Does that sports team belong to Sacile? S: No	INTENT =Verification 2 entities, subject is indirect ENT =[Q653772 (Pittsburgh Pirates), Q53190 (Sacile)], REL =[P17 (country)], TYP =[Q15617994 (designation admin. territorial entity)], TRIPLE =[(Q653772,P17,Q53190)], GOLD =[False]	AS: [is_in, Q53190, find, Q653772, P17] SP: ASK {wd:Q653772 wdt:P17 wdt:Q53190.}

Table 1: Example conversations from SPICE. The left column shows the sequence of user (U) and system (S) utterances. The middle column shows the annotations provided in CSQA. Blue boxes on the right show the sequence of actions (AS) and corresponding SPARQL semantic parses (SP).

Nb. instances	197K
Nb. entities	12.8M
Nb. relations	2738
Nb. types	3064
Avg. turn length	9.5
Avg. entities per conversation	7.6
Avg. types per conversations	6.5
Avg. neighbourhood per turn	181.4 triples
Logical Reasoning, Quantitative Reasoning, Comparative Reasoning, Quantitative Reasoning Count, Comparative Reasoning Count, Verification, Simple Question	
Clarification, Coreference, Ellipsis	

Table 2: Statistics of SPICE dataset (top); general question types (middle); linguistic phenomena (bottom).

Conversations contain sequences of mixed-initiative interactions where the system requests clarification on ambiguous questions. Conversations also include discourse phenomena such as coreference (e.g., *Which sports team was the champion of that tournament?* in Table 1) and ellipsis (e.g., *And what about 1910 World Series?* as a follow up question for *How many sports teams participated in that tournament?*).

In total, there are 10 question types and 47 question sub-types. In Table 2, we only list question types but provide all question sub-types in Table 10 in the Appendix A.

2.1 Question Semantics Described by Actions

Saha et al. (2018) envisaged CSQA as a benchmark for retrieval-based conversational question answering (Bordes et al., 2015; Dong et al., 2015; Jain, 2016; Lan et al., 2021). These methods em-

bed natural language questions and KG triples into high dimensional spaces and rely on neural reasoning modules to match the question to candidate answers. Hence, questions do not have associated logical forms describing their meaning, only gold answers are available.

Our success in creating semantic parse annotations is partly due to the fact that CSQA provides useful KG information. Each interaction (i.e., user and system turn) comes with annotations about KG entities, types, and relation symbols as well as some information about the triple patterns involved in the question (illustrated in Table 1 with ENT, REL, TYP, and TRIPLE fields). It also provides information pertaining to question types and sub-types (see INTENT in Table 1).

Taking advantage of these annotations, follow-on work (Guo et al., 2018) defined a semantic parsing task over CSQA, modeling the meaning of questions as a sequence of actions. The set of actions encompasses find (or find_rev when the entity is in object position) to retrieve sets of entities in a subject (object) position, as well as actions operating on sets of entities (e.g., filter_type). For instance, the question in turn \mathcal{T}_1 in Table 1 would be parsed to [filter_type, find_rev, Q650855, P1923, Q500834], meaning “find the set of entities that are in relationship *participating team* with *Detroit Tigers* and then filter those that are of type *tournament*”. A breath-first search algorithm generates action-grammar annotations for each question and a sequence of grammar-actions is considered correct if upon exe-

	ATIS	SPaC	CoSQL	SPICE
Nb. Instances	1,658	4,298	3,007	197K
Avg. turn length	7.0	3.0	5.2	9.5
Domain	Single	Multi	Multi	Wikidata
Logical form	SQL	SQL	SQL	SPARQL
Database type	Rel	Rel	Rel	KG

Table 3: Datasets for conversational semantic parsing (Rel: relational database; KG: knowledge graph).

cution it returns the gold answer. Subsequent work expanded the action-grammar (Shen et al., 2019; Kacupaj et al., 2021; Marion et al., 2021), greatly improving its coverage (i.e., the number of successfully annotated questions).

2.2 From Actions to SPARQL Queries

In this work, we take a step further and map CSQA natural language questions into SPARQL queries. Being a standard query language for KGs, SPARQL affords generality. The dataset can be further extended with new question types, and existing models adapted to these, without developing a new grammar and its corresponding execution engine. In addition, semantic parsing approaches could be compared on an equal footing as different grammars might have different coverage and each action encapsulate different degrees of complexity.

We first analysed how intent is expressed in question types and sub-types as well as the action-grammar derivations thereof. We then manually defined SPARQL templates for each question sub-type. A SPARQL template is a query with unspecified triple patterns in the WHERE clause. For instance, the SPARQL template for the question in turn \mathcal{T}_1 is {SELECT ?x WHERE ?x RELATION ENTITY. ?x wdt:P31 TYPE.}. We then modified the tool provided in Kacupaj et al. (2021) to automatically instantiate the SPARQL templates, providing annotations for the entire dataset. For each question, the tool searches for combinations over the KG symbols provided in the CSQA dataset to derive sequences of actions (e.g., [filter_type, find_rev, Q650855, P1923, Q500834] for question in \mathcal{T}_1 in Table 1). In our case, the tool fills missing slots in SPARQL templates (e.g., SELECT ?x WHERE {?x wdt:P1923 wd:Q650855. ?x wdt:P31 wd:Q500834.}).

We assessed the correctness of instantiated SPARQL queries by executing them and comparing results to gold answers. To be able to execute SPARQLs, we imported the Wikidata snapshot provided by Saha et al. (2018) into a KG in a SPARQL server (see Appendix B for more details). For some

questions, the annotation procedure did not produce a SPARQL parse that recovers the gold answer. In such cases, which are a minor percentage of the entire dataset, we redefine the answer whenever this does not affect the conversation flow or truncate the conversation at the question that could not be annotated.

Table 2 shows various statistics for our SPICE dataset. Compared to related conversational datasets such as ATIS (Suhr et al., 2018), SPaC (Yu et al., 2019b), and CoSQL (Yu et al., 2019a) (see Table 3), SPICE contains a sizeable number of training instances, conversations are longer, and the semantic parsing task is real-scale.

3 The Semantic Parsing Task

We consider the semantic parsing task over a sequence of dialogue turns $d = (d_1, d_2, \dots, d_{|d|})$, where turn d_t corresponds to a user-system interaction with user question x_t and system answer a_t . Each turn has a conversation context c_t made of interactions d_i such that $i < t$. Given interaction d_t with context c_t and user question $x_t = (x_{t1}, x_{t2}, \dots, x_{t|x_t|})$, our goal is to predict a SPARQL query (logical form) $y_t = (y_{t1}, y_{t2}, \dots, y_{t|y_t|})$ that represents the intent of x_t and, upon execution over knowledge-graph \mathcal{K} , yields denotation (answer) a_t . y_t is a sequence over a target vocabulary $\mathcal{V} = \mathcal{V}_f \cup \mathcal{V}_K$ where \mathcal{V}_f is a fixed vocabulary containing SPARQL keywords (e.g., SELECT) and additional special tokens (e.g., beginning of sequence token, BOS) and \mathcal{V}_K contains all knowledge-graph symbols (e.g., entity IDs such as Q76 for *Barack Obama*).

We propose two approaches for this semantic parsing task which establish strong baseline performance and highlight various challenges. These differ in the way they handle large KG vocabularies and how they generate logical forms. Figure 1 provides a graphical overview of the two models discussed below.

3.1 Parsing with a Single Decoder and Knowledge Sub-graphs

Our first semantic parsing model is parametrised by an encoder-decoder Transformer neural network (Vaswani et al., 2017). It is an adaptation of the architecture proposed in Gu et al. (2021) for standard semantic parsing. This approach relies on *dynamic vocabularies* to deal with the large target vocabulary inherent in semantic parsing over large scale

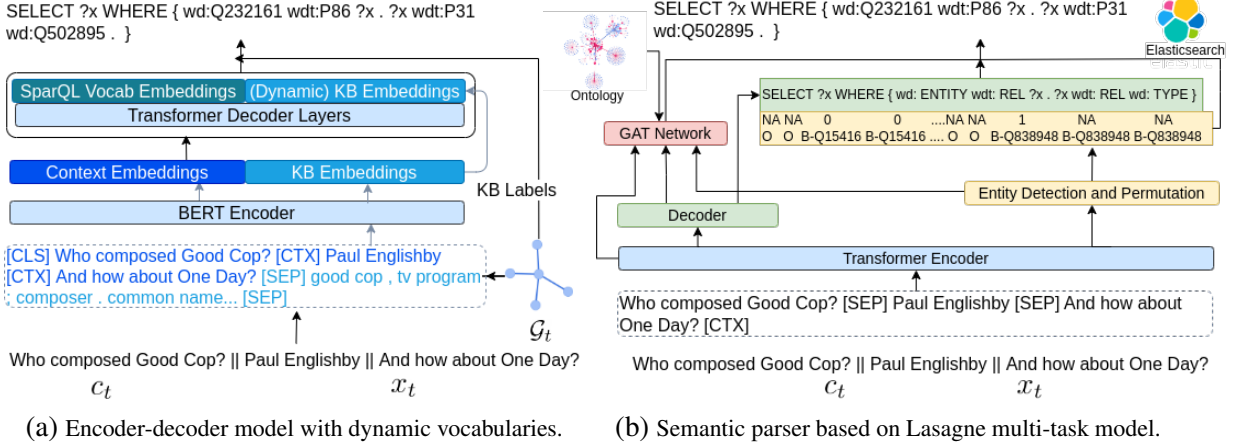


Figure 1: Two modeling approaches to conversational semantic parsing.

KGs containing thousands of relations and types, and millions of entities.

Dynamic Vocabulary Since the KG vocabulary \mathcal{V}_K can be extremely large, we parse question x_t with a smaller vocabulary $\mathcal{V}_t \subseteq \mathcal{V}_K$ which only contains KG symbols related to x_t . Following previous work (Gu et al., 2021; Marion et al., 2021), we assume the symbols related to x_t are those appearing in sub-graph \mathcal{G}_t of knowledge-graph \mathcal{K} , $\mathcal{G}_t \subseteq \mathcal{K}$. Given question x_t and its context c_t , we identify KG entities $\mathcal{E}_t = \{e_{t1}, e_{t2}, \dots, e_{t|\mathcal{E}_t|}\}$ which correspond to mentions in x_t and c_t . We then obtain \mathcal{G}_t by taking the one-hop neighbourhood for each entity $e_{ti} \in \mathcal{E}_t$. In other words, we include all KG triples (s, r, o) where the entity appears in the subject ($s = e_{ti}$) or object position ($o = e_{ti}$). When e_{ti} is a subject, we include triple (e_{ti}, r, τ_o) where τ_o is the type of entity o ; analogously, when e_{ti} appears in the object position, we add (τ_s, r, e_{ti}) . For entities e_{ti} we include their types $\tau_{e_{ti}}$. When e_{ti} is a general entity (e.g., a type such as *tournament*) we add those relations from \mathcal{K} that have instances of type e_{ti} as their subject (object). The final vocabulary \mathcal{V}_t contains all entities in \mathcal{E}_t and all relations r and types (τ_o, τ_s , and $\tau_{e_{ti}}$) found in the set of triples in \mathcal{G}_t .

Note that context c_t is defined as a window over the conversation so far. Following previous work (Marion et al., 2021; Kacupaj et al., 2021), we set the conversation context to the previous user-system interaction $c_t = \{d_{t-1}\}$.

Encoder-Decoder Model Our encoder is a BERT (Devlin et al., 2019) model fine-tuned on our semantic parsing task. The decoder is a randomly initialised Transformer network (Vaswani

et al., 2017). To account for the difference in initialisation between the encoder and decoder networks, we follow the training scheme proposed in Liu and Lapata (2019). We provide details in Appendix C.

The input to our semantic parser is a tuple $(x_t, c_t, \mathcal{G}_t)$ consisting of natural language question x_t , its context c_t , and subgraph \mathcal{G}_t which we adapt to BERT’s input format as follows (Gu et al., 2021). We concatenate the sequence of natural language questions and answers appearing in c_t and x_t , using the special token [CTX] as a delimiter and prepend the [CLS] token in the beginning of the sequence. Special token [SEP] denotes the end of sequence followed by the linearised KG sub-graph \mathcal{G}_t . The linearisation procedure goes over entities in \mathcal{G}_t , enumerating their types and relations. Importantly, we denote entities by their label rather than their KG identifiers. The order of entities in \mathcal{G}_t is random. Figure 1(a) shows an example of the input to our BERT-based encoder.

More formally, the encoder takes token sequences $x'_t = [\text{CLS}]x'_{t_{\text{text}}}[\text{SEP}]x'_{t_{\text{graph}}}[\text{SEP}]$ as input where $x'_{t_{\text{text}}}$ is the natural language subsequence $(x_1^{t-1}, \dots, x_{|x_t|-1}^{t-1}, a_1^{t-1}, \dots, a_{|a_t|-1}^{t-1}, x_1^t, \dots, x_{|x_t|}^t)$ and $x'_{t_{\text{graph}}} = (g_1^t, \dots, g_{|\mathcal{G}_t|}^t)$ is the sequence of knowledge-graph symbols from the linearised graph \mathcal{G}_t . Note that these knowledge-graph symbols constitute the target dynamic vocabulary \mathcal{V}_t and $|\mathcal{G}_t|$ represents the number of KG symbols which is equal to the size of the target vocabulary $|\mathcal{V}_t|$. The encoder maps input sequences x'_t into sequences of continuous representations $\mathbf{z}_t = (\mathbf{z}_{t1}, \dots, \mathbf{z}_{t|x_t|})$, and the decoder then generates the target SPARQL parse $y_t = (y_{t1}, \dots, y_{t|y_t|})$ token-by-token autoregressively, hence modelling the conditional probability: $p(y_{t1}, \dots, y_{t|y_t|} | x'_t)$.

The linearised graph \mathcal{G}_t can exceed BERT’s maximum number of input positions (which is 512). To avoid throwing away useful information, we adopt a solution similar to Gu et al. (2021). For question x_t with \mathcal{G}_t containing k entities, we create k input sequences x_t^1, \dots, x_t^k . These k sequences share the natural language subsequence but have different KG symbol subsequences. Given an input sequence x_t^1, \dots, x_t^k , we obtain contextualised representations as $\mathbf{z}_t^1, \dots, \mathbf{z}_t^k = \text{BERT}(x_t^1, \dots, x_t^k)$.

The model further splits the sequence of continuous representations \mathbf{z}_t^j into textual representations $\mathbf{z}_{t_{\text{text}}}^j$ and knowledge-graph symbols $\mathbf{z}_{t_{\text{graph}}}^j$ both of which are contextualised. We then average representations $\mathbf{z}_{t_{\text{text}}} = \text{AVG}(\mathbf{z}_{t_{\text{text}}}^j)$ and feed them as input to the decoder (see Figure 1(a)). From representations $\mathbf{z}_{t_{\text{graph}}}^j$, we derive the embeddings for the elements in the target dynamic vocabulary \mathcal{V}_t . Recall that the decoder parses input questions x_t using target vocabulary $\mathcal{V} = \mathcal{V}_f \cup \mathcal{V}_t$ which consists of a set of fixed (\mathcal{V}_f) and dynamic (\mathcal{V}_t) target tokens. The decoder then predicts the probability of each SPARQL token y_{ti} as follows:

$$p(y_{ti} | y_{t < i}, x_t^1, \dots, x_t^k) = \text{softmax}(\mathbf{W}_o \mathbf{h}_i^L)$$

where \mathbf{h}_i^L is the decoder top layer hidden representation at time step i . $\mathbf{W}_o \in \mathbb{R}^{|\mathcal{V}_f \cup \mathcal{V}_t|}$ is the output embedding matrix with $\mathbf{W}_o = [\mathbf{W}_f; \mathbf{W}_t]$, where $[\cdot]$ denotes matrix concatenation, and \mathbf{W}_t is derived from the encoder representations $\mathbf{z}_{t_{\text{graph}}}^j$.

3.2 Parsing with Multiple Decoders and an Ontology Classifier

Our second model is an adaptation of the Lasagne architecture proposed in Kacupaj et al. (2021). Lasagne generates logical forms following a multi-stage approach where a backbone sketch is first predicted and then fleshed out. Their sketch is a sequence of actions from a custom grammar which we modify to predict SPARQL queries.

Action Prediction Lasagne employs an encoder-decoder model based on Transformers (Vaswani et al., 2017) to convert a user question x_t in a conversation into a logical form template. The input to the encoder is the the conversation context $c_t = \{d_{t-1}\}$ and user question x_t . Utterances are separated via [SEP] tokens, while the special context token [CTX] denotes the end of sequence (see Figure 1b). The input sequence is encoded via

multi-head attention (Vaswani et al., 2017) to output contextualized representations which are then fed to the decoder to predict logical forms templates (without grounding to KG elements) token-by-token.

While Kacupaj et al. (2021)’s model predicts a sequence of actions valid in their custom grammar (e.g., filter_type, find_rev) our decoder predicts SPARQL queries. Our decoder predicts SPARQLs with place-holders for KG symbols. For instance, for the WHERE clause of turn \mathcal{T}_1 in Table 1, it predicts {ENTITY RELATION ?x. ?x wdt:P31 TYPE.} instead of {wd:Q5582479 wdt:P161 ?x. ?x wdt:P31 wd:Q502895.})

Entity Recognition and Linking An entity recognition module detects entities in the input and links them to the KG (Shen et al., 2019; Kacupaj et al., 2021). Initially, entity spans are identified using an LSTM which performs BIO sequence labelling.³ Entity spans are subsequently linked to KG entities via an inverted index (created using Elasticsearch⁴) which maps entity labels to entity IDs. Once identified, the entities are further filtered and reordered so that they match their order of appearance in the SPARQL (see Figure 1(b)).

Predicting Types and Relations Finally, an ontology graph including types and relations appearing in SPICE’s KG is constructed.⁵ The graph is encoded with a Graph Attention Network (GAT; Velickovic et al. 2018) and the prediction of type and relation fillers for the SPARQL template is modeled as a classification task over graph nodes, given the current conversational context and the decoder hidden state.

Learning All modules outlined above are trained in a multi-task manner, optimizing the weighted average of the following individual losses:

$$L = \lambda_1 L^F + \lambda_2 L^G + \lambda_3 L^R + \lambda_4 L^O$$

where L^F is the loss of the logical form template decoder, L^G is the type and relation prediction loss using the GAT network, L^R is the entity recognition loss, and L^O the entity reordering loss (and weights $\lambda_{1:4}$ are learned during training). We refer

³BIO labels for training are obtained by performing string matching between entity annotations and user utterances.

⁴<https://www.elastic.co/>

⁵Note that for a semantic parsing system in production over the full Wikidata KG this graph could be significantly bigger.

the interested reader to [Kacupaj et al. \(2021\)](#) for details on model initialisation and training.

4 Results

We examine how the two models just described fare on different question types and subtypes. We report results on SPICE i.i.d train/valid/test splits (containing 152,391/16,813/27,797 conversations, respectively) but also create new splits that assess out-of-distribution generalization. In all cases, following previous work ([Saha et al., 2018](#); [Kacupaj et al., 2021](#)), we use execution-based automatic metrics. *F1-score* evaluates question parses that return a set of entities, while *Accuracy* is used for question parses that evaluate to True/False or return a numerical value.⁶ In addition, we report *Exact Match* (EM) against the gold SPARQL parse.

4.1 Performance per Question Type

Table 4 shows our results on the SPICE i.i.d test split. We evaluate the BERT-based semantic parser with dynamic vocabularies (BertSP) and the semantic parser based on the Lasagne system (LasagneSP) discussed in Section 3. The BertSP_G variant has access to oracle entities, types, and coreference annotations which allows us to disassociate the complexity of the SPARQL generation task from the problem of grounding and disambiguating entities to KG symbols. The other two variants do not have access to these oracle annotations. To ground named entity mentions present in user questions to KG entities, BertSP_S uses a vanilla string-matching based Named Entity Linking (NEL) algorithm similar to that used by [Marion et al. \(2021\)](#) while BertSP_A relies on an off-the-shelf Name Entity Recognition (NER) system and an inverted index for Named Entity Linking (NEL).⁷ None of these variants has access to oracle coreference entities and have to resolve them using the conversation context c_t . Both use a string-matching based type (i.e., general entities) linking approach. It is not straightforward to perform oracle analysis for LasagneSP without compromising the model structure which predicts entities, their types, and relations in multiple stages.

Performance Assessment with Exact Match.

We observe that execution based metrics (F1-score and Accuracy) are generally higher than EM. This is because in some cases the SPARQL parse may be

incorrect and still yield some results. For instance, a parse requiring the UNION of two graph patterns may yield a partially correct answer by only including one graph pattern; similarly, a parse can evaluate to False and agree with the gold answer just because it included a wrong relation symbol.

Importance of Entity Grounding. Not surprisingly, the model with access to oracle information (variant BertSP_G) obtains the best performance. Results improve not only for questions with entities referring to previous context but also indirectly for other types of questions. Since entities are correctly grounded in previous conversation turns c_t , the model operates with more accurate graphs \mathcal{G}_t and richer dynamic vocabularies \mathcal{V}_t (c.f. Section 3.1). Both BertSP_S and BertSP_A have to solve coreference using a limited conversation context and thus performance decreases. These also have to solve named (*Detroit Tigers*) and non-named general (*tournament*) entity mentions. As discussed above, they differ in the approach to NER/NEL (vanilla string matching for BertSP_S and an off-the-shelf NER/NEL system for BertSP_A). The performance of BertSP_A is worse than BertSP_S. We observe in the data that this is related to NER errors (specially compound named entities such as *President of the Czech Republic*) and disambiguation during NEL (e.g., *Saint Barbara* the painting versus the Saint).⁸ While previous work has relied on vanilla string matching and custom approaches to NER/NEL with great success, we encourage research to rely on NER/NEL methods with application in realistic scenarios (i.e., that work for different surface forms of entity mentions and on large scale KGs where latency and disambiguation performance are crucial).

Differences in Modelling Approach. BertSP_S is closer to LasagneSP in terms of how they solve NER/NEL with a task specific approach; however, their semantic parsing approach is conceptually different (c.f. Section 3). LasagneSP outperforms BertSP_S in Comparative, Quantitative, and Comparative-Count questions. These encompass many question sub-types with general non-named entities which are very common in training and test splits. LasagneSP has access to all types and relations encoded with the graph network. In contrast, BertSP_S relies on types that firstly need to

⁶Following previous work we report Micro *F1-score*.

⁷We use AllenNLP’s NER model and Elasticsearch.

⁸Note that for better recall, for each NER entity we include the top-K linked KG entities.

Question Type	BertSP _G		BertSP _S		BertSP _A		LasagneSP	
	F1	EM	F1	EM	F1	EM	F1	EM
Clarification	84.89	82.53	80.21	77.69	83.91	76.58	86.29	73.41
Logical Reasoning (All)	90.61	82.90	85.55	66.89	22.74	28.61	88.80	57.41
Quantitative Reasoning (All)	94.42	88.55	82.95	66.40	76.20	59.01	94.90	91.47
Comparative Reasoning (All)	96.23	87.39	90.44	73.80	69.56	39.37	94.20	85.05
Simple Question (Coreferenced)	88.96	86.53	83.19	69.87	76.51	58.83	84.73	60.90
Simple Question (Direct)	91.81	91.59	87.13	80.69	71.43	58.71	87.21	66.88
Simple Question (Ellipsis)	79.51	89.71	72.50	71.67	58.14	50.90	74.35	61.53
	AC	EM	AC	EM	AC	EM	AC	EM
Verification (Boolean)	90.10	77.24	79.72	62.62	37.16	24.90	34.89	26.72
Quantitative Reasoning (Count)	87.91	84.97	76.88	73.20	50.86	48.44	60.51	56.15
Comparative Reasoning (Count)	90.05	85.99	73.18	66.79	43.48	40.67	89.09	83.69
Overall	81.50	85.74	81.18	70.96	59.00	48.60	79.50	66.32

Table 4: Results on SPICE i.i.d test split. F1-Score (F1), accuracy (AC), and exact match (EM).

be present in the entity neighbourhood sub-graph and then be preserved after input truncation to fit BertSP_S input size. An advantage of BertSP_S architecture over LasagneSP, is that it allows for easier adaptation to new types and relations by relying on dynamic vocabularies. Architectures as LasagneSP based on type and relation classifiers need to be retrained to accommodate for these new elements. In Simple questions, where each question involves fewer but more diverse types, BertSP_S input text and KG symbols contextualisation helps in predicting more accurate types and thus BertSP_S achieves better performance. LasagneSP performs poorly on Verification, Logical, and Quantitative-Count (this last encompasses logical operators). This is explained by LasagneSP’s modelling limitations, i.e., not been able to point to the same input entity more than once.

Observed Errors in Predicted SPARQLs. A manual inspection into incorrect SPARQL predictions reveals as common errors the prediction of wrong entities and relations, failure to enumerate all required entities (for questions with multiple entities), and wrong argument order (i.e., entities and variables are correct but placed in the incorrect subject/object positions of graph patterns); to a lesser extent SPARQL predictions with incorrect question intent and ill-formed formulas are also seen.

4.2 Linguistic Analysis

Table 5 shows model performance across-different question sub-types aggregated for specific phenomena. These include coreference, ellipsis, and multiple entities. We distinguish between cases where coreference can be resolved in the previous turn (d_{t-1}) and further back in the conversation history (d_{t-i} , where $i > 1$). In addition, some question sub-types contain plural mentions, i.e., they

Phenomena	BertSP _G	BertSP _S	BertSP _A	LasagneSP
Coreference ₌₋₁	81.40	70.65	49.39	43.65
Coreference _{<-1}	67.82	0	0	0
Ellipsis	75.93	54.33	26.39	46.54
Multiple Entities	83.37	65.40	41.64	66.52

Table 5: Results (exact match) on SPICE i.i.d test; questions grouped according to linguistic phenomena.

are linked to multiple entities which the semantic parser must enumerate in order to build the correct parse. Ellipsis can be often resolved within the previous interaction (d_{t-1}). Questions with multiple entities bring further disambiguation challenges. In Appendix A we provide the list of question sub-types we aggregate in each phenomena.

As Table 5 shows, an oracle model that has access to gold annotations is superior to models which rely on automatic entity and type linking. BertSP_S performs better than LasagneSP at handling coreference that can be resolved within the immediate context (i.e., $c_t = \{d_{t-1}\}$). Due to the fact that LasagneSP predicts entity positions in SPARQL, it is particularly bad at resolving mentions to multiple entities in the previous context or even the same entity appearing multiple times in the output parse (as is the case with Verification questions). Perhaps unsurprisingly, neither BertSP_S nor LasagneSP can resolve mentions to non-immediately preceding utterances. BertSP_S performs better than LasagneSP in questions with Ellipsis, we conjecture that its input context and KG symbols contextualisation help in grounding ellided relation mentions. Ellipsis and Multiple Entities improve by a large margin with access to gold annotations (see BertSP_G in Table 5).

4.3 Generalisation

We further evaluate the models’ ability to generalize by creating “query-based” splits (Finegan-Dollak et al., 2018), i.e., splits with query patterns

Unseen Combinations	Instances Train/valid/test	BertSP _S EM	LasagneSP EM
COUNTLOGIC	153,562/14,262/29,177	0.94	0
UNIONMULTI	157,331/14,426/25,244	19.73	16.89
VERIFY3	154,027/13,869/29,105	0	0

Table 6: Model performance on SPICE non-i.i.d splits.

seen only at test time. Our compositional splits are shown in Table 6, and include: (a) question sub-types that involve a count operation over a union operator (COUNTLOGIC; individual operators are seen at training time but not the combination thereof); (b) question sub-types that involve a union operator over two graph patterns with different relations (UNIONMULTI; the union of two graph patterns with the *same* relation is seen during training); and (c) verification questions with three entities (VERIFY3; questions with 2 entities only are seen during training).

As shown in Table 6, both BertSP_S and LasagneSP perform poorly across different splits. This suggests that they successfully memorise SPARQL parses seen during training but are unable to generalize to similar user questions with some language variation. While models, in some cases, grasp the overall SPARQL structure for the unseen user question (e.g., *Which watercourses are located in the neighbourhood of Bremen or are the tributaries of Ob? in UNIONMULTI*), they ignore the different details in the question and produce the final SPARQL systematically in the same way as done for similar seen questions (e.g., *Which people are the creators of The Theory of Everything or Ten Minutes to Live?*). In UNIONMULTI models produce an appropriate SPARQL template but systematically copy the *same* relation in both graph patterns. BertSP_S performs slightly better than LasagneSP as BertSP_S’s contextualised KG embeddings may help to chose the different relations in the graph patterns. Similar trend is observed for COUNTLOGIC and VERIFY3. In Appendix D we provide examples of unseen user questions, SPARQLs, and common errors.

5 Related Work

Much previous work on semantic parsing focuses on mapping stand-alone utterances to logical forms. Relatively few datasets have been constructed for *conversational* semantic-parsing (Suhr et al., 2018; Dahl et al., 1994; Yu et al., 2019b,a) partly due to the difficulty of eliciting annotations in an interac-

tive context. As a result, existing benchmarks are either single-domain or small-scale (see the comparison in Table 3). For instance, ATIS (Suhr et al., 2018) is a conversational dataset, it contains utterances paired with SQL queries pertaining to a US flight booking task and exemplifies several challenging long-range discourse phenomena (Jain and Lapata, 2021); however, it is restricted to a single domain with a simple database schema. SPaC (Yu et al., 2019b) and CoSQL (Yu et al., 2019a) cover multiple domains and include multi-turn user and system interactions. These datasets present cross-domain challenges in mapping natural language queries onto SQL, but the conversation length is fairly short and the databases relatively small-scale.

Large KGs, like Wikidata (Vrandečić, 2012), are becoming an increasing valuable source of information. Various question-answering datasets have been recently released with complex questions (and their semantic parses) grounded to knowledge bases (Dubey et al. 2019; Talmor and Berant 2018, *inter alia*). However, these do not cover conversational queries. The CSQA dataset introduced in Saha et al. (2018) is conversational, covering a wide range of linguistic phenomena (e.g., ellipsis, coreference) but frames the QA task as an information retrieval problem. Follow-on work (Marion et al., 2021; Kacupaj et al., 2021; Shen et al., 2019) has devised a semantic parsing task for this dataset using a hand-crafted grammar to automatically obtain annotations by searching and combining rules in a breadth first search manner. However, as discussed in Section 2, these are not directly executable with a real KG engine, such as Blazegraph⁹, and are further restricted by the underlying grammar which would have to be (perpetually) tweaked to improve coverage for new types of questions and conversations.

6 Conclusion

In this work we introduce SPICE, a conversational semantic parsing dataset over knowledge graphs. Our dataset contains SPARQL annotations which are executable on a real KG engine and requires handling complex questions, type, relation, and entity linking on large scale. Moreover, it showcases multiple linguistic phenomena such as coreference and ellipsis. We establish two strong baselines for the semantic parsing task and present detailed anal-

⁹We use Blazegraph (<https://blazegraph.com/>) for hosting our data and executing SPARQLs.

ysis stratifying performance by question type and linguistic phenomena. We also study generalization to unseen question intents and create multiple dataset splits with different query patterns. We hope our dataset will serve as useful testbed for the development of conversational semantic parsers.

References

- Antoine Bordes, Nicolas Usunier, Sumit Chopra, and Jason Weston. 2015. [Large-scale simple question answering with memory networks](#). *CoRR*, abs/1506.02075.
- Deborah A. Dahl, Madeleine Bates, Michael Brown, William Fisher, Kate Hunicke-Smith, David Pallett, Christine Pao, Alexander Rudnicky, and Elizabeth Shriberg. 1994. [Expanding the scope of the ATIS task: The ATIS-3 corpus](#). In *Human Language Technology: Proceedings of a Workshop held at Plainsboro, New Jersey, March 8-11, 1994*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Li Dong, Furu Wei, Ming Zhou, and Ke Xu. 2015. [Question answering over Freebase with multi-column convolutional neural networks](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 260–269, Beijing, China. Association for Computational Linguistics.
- Mohnish Dubey, Debayan Banerjee, Abdelrahman Abdelkawi, and Jens Lehmann. 2019. Lc-quad 2.0: A large dataset for complex question answering over wikidata and dbpedia. In *The Semantic Web – ISWC 2019*, pages 69–78, Cham. Springer International Publishing.
- Catherine Finegan-Dollak, Jonathan K. Kummerfeld, Li Zhang, Karthik Ramanathan, Sesh Sadasivam, Rui Zhang, and Dragomir Radev. 2018. [Improving text-to-SQL evaluation methodology](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 351–360, Melbourne, Australia. Association for Computational Linguistics.
- Yu Gu, Sue Kase, Michelle Vanni, Brian Sadler, Percy Liang, Xifeng Yan, and Yu Su. 2021. [Beyond i.i.d.: Three levels of generalization for question answering on knowledge bases](#). In *Proceedings of the Web Conference 2021, WWW '21*, page 3477–3488, New York, NY, USA. Association for Computing Machinery.
- Daya Guo, Duyu Tang, Nan Duan, Ming Zhou, and Jian Yin. 2018. [Dialog-to-action: Conversational question answering over a large-scale knowledge base](#). In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Parag Jain and Mirella Lapata. 2021. [Memory-Based Semantic Parsing](#). *Transactions of the Association for Computational Linguistics*, 9:1197–1212.
- Sarthak Jain. 2016. [Question answering over knowledge base using factual memory networks](#). In *Proceedings of the NAACL Student Research Workshop*, pages 109–115, San Diego, California. Association for Computational Linguistics.
- Dietmar Jannach, Ahtsham Manzoor, Wanling Cai, and Li Chen. 2021. [A survey on conversational recommender systems](#). *ACM Comput. Surv.*, 54(5).
- Endri Kacupaj, Joan Plepi, Kuldeep Singh, Harsh Thakkar, Jens Lehmann, and Maria Maleshkova. 2021. [Conversational question answering over knowledge graphs with transformer and graph attention networks](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 850–862, Online. Association for Computational Linguistics.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Yunshi Lan, Gaole He, Jinhao Jiang, Jing Jiang, Wayne Xin Zhao, and Ji-Rong Wen. 2021. [A survey on complex knowledge base question answering: Methods, challenges and solutions](#). In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 4483–4491. International Joint Conferences on Artificial Intelligence Organization. Survey Track.
- Yang Liu and Mirella Lapata. 2019. [Text summarization with pretrained encoders](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3730–3740, Hong Kong, China. Association for Computational Linguistics.
- Pierre Marion, Pawel Nowak, and Francesco Piccinno. 2021. [Structured context and high-coverage grammar for conversational question answering over knowledge graphs](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8813–8829, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. [Pytorch: An imperative style, high-performance deep learning library](#). In *Advances in Neural Information Processing Systems*, volume 32, pages 8026–8037. Curran Associates, Inc.
- Filip Radlinski and Nick Craswell. 2017. [A theoretical framework for conversational search](#). In *Proceedings of the 2017 Conference on Conference Human Information Interaction and Retrieval, CHIIR '17*, page 117–126, New York, NY, USA. Association for Computing Machinery.
- Siva Reddy, Mirella Lapata, and Mark Steedman. 2014. [Large-scale semantic parsing without question-answer pairs](#). *Transactions of the Association for Computational Linguistics*, 2:377–392.
- Amrita Saha, Vardaan Pahuja, Mitesh Khapra, Karthik Sankaranarayanan, and Sarath Chandar. 2018. [Complex sequential question answering: Towards learning to converse over linked question answer pairs with a knowledge graph](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18)*, pages 705–713, New Orleans, Louisiana, USA. AAAI Press.
- Tao Shen, Xiubo Geng, Tao Qin, Daya Guo, Duyu Tang, Nan Duan, Guodong Long, and Daxin Jiang. 2019. [Multi-task learning for conversational question answering over a large-scale knowledge base](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2442–2451, Hong Kong, China. Association for Computational Linguistics.
- Alane Suhr, Srinivasan Iyer, and Yoav Artzi. 2018. [Learning to map context-dependent sentences to executable formal queries](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2238–2249, New Orleans, Louisiana. Association for Computational Linguistics.
- Alon Talmor and Jonathan Berant. 2018. [The web as a knowledge-base for answering complex questions](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 641–651, New Orleans, Louisiana. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. [Graph attention networks](#). In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.
- Denny Vrandečić. 2012. [Wikidata: A new platform for collaborative data collection](#). In *Proceedings of the 21st International Conference on World Wide Web, WWW '12 Companion*, page 1063–1064, New York, NY, USA. Association for Computing Machinery.
- Wen-tau Yih, Matthew Richardson, Chris Meek, Ming-Wei Chang, and Jina Suh. 2016. [The value of semantic parse labeling for knowledge base question answering](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 201–206, Berlin, Germany. Association for Computational Linguistics.
- Tao Yu, Rui Zhang, Heyang Er, Suyi Li, Eric Xue, Bo Pang, Xi Victoria Lin, Yi Chern Tan, Tianze Shi, Zihan Li, Youxuan Jiang, Michihiro Yasunaga, Sungrok Shim, Tao Chen, Alexander Fabbri, Zifan Li, Luyao Chen, Yuwen Zhang, Shreya Dixit, Vincent Zhang, Caiming Xiong, Richard Socher, Walter Lasecki, and Dragomir Radev. 2019a. [CoSQL: A conversational text-to-SQL challenge towards cross-domain natural language interfaces to databases](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1962–1979, Hong Kong, China. Association for Computational Linguistics.
- Tao Yu, Rui Zhang, Michihiro Yasunaga, Yi Chern Tan, Xi Victoria Lin, Suyi Li, Heyang Er, Irene Li, Bo Pang, Tao Chen, Emily Ji, Shreya Dixit, David Proctor, Sungrok Shim, Jonathan Kraft, Vincent Zhang, Caiming Xiong, Richard Socher, and Dragomir Radev. 2019b. [SPaRC: Cross-domain semantic parsing in context](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4511–4523, Florence, Italy. Association for Computational Linguistics.
- Hamed Zamani, Johanne R. Trippas, Jeff Dalton, and Filip Radlinski. 2022. Conversational information seeking. <https://arxiv.org/abs/2201.08808>.

A The SPICE Dataset: Question Types and Sub-Types

Table 10 provides the list of question types and sub-types in SPICE. For each question sub-type we provide an example user question. For cases involving ellipsis and coreference, we include the conversation context (in grey colour).

Table 7 provides the list of question sub-types grouped per linguistic phenomena. Coreference ($=_{-1}$ and $<_{-1}$), Ellipsis, and Multiple Entities.

Coreference
More/Less Mult. entity type (Coreference) # More/Less Single entity type (Coreference) # Single Entity (Coreference) # 2 entities, one direct and one indirect, object is indirect # 2 entities, one direct and one indirect, subject is indirect # 3 entities, 2 direct, 2(direct) are query entities, subject is corefered # one entity, multiple entities (as object) corefered # Count Logical operators (Coreference) # Count Single entity type (Coreference) # Count over More/Less Mult. entity type (Coreference) # Count over More/Less Single entity type (Coreference)
Ellipsis
Difference Single Relation (Ellipsis) # Intersection Single Relation (Ellipsis) # Union Single Relation (Ellipsis) # More/Less Mult. entity type (Ellipsis) # More/Less Single entity type (Ellipsis) # object parent is changed, subject and predicate remain same # Incomplete count-based ques # Count over More/Less Mult. entity type (Ellipsis) # Count over More/Less Single entity type (Ellipsis)
Multiple Entities
Difference Multiple Relation # Intersection Multiple Relation # Union Multiple Relation # Atleast/ Atmost/ Approx. the same/Equal Mult. entity type # Min/Max Mult. entity type # More/Less Mult. entity type # More/Less Mult. entity type (Ellipsis) # More/Less Mult. entity type (Coreference) # Mult. Entity (Simple Question Direct and Coreference) # one entity, multiple entities (as object) coreferred # Count over Atleast/ Atmost/ Approx. the same / Equal Mult. entity type # Count Mult. entity type # Count over More/Less Mult. entity type # Count over More/Less Mult. entity type (Ellipsis) # Count over More/Less Mult. entity type (Coreference)

Table 7: Unseen question sub-types in SPICE non-i.i.d splits.

Table 8 provides the list of unseen question sub-types for each of the non-i.i.d splits.

B Creating a Knowledge Graph from the CSQA Data

Deploying a full copy of Wikidata locally as a standalone service requires huge resources in addition to cluster dependent tweaking to obtain fast query processing and high-performance.¹⁰ To enable easier deployment and fast access for research purposes we created a smaller graph from the CSQA data files. We mapped the contents of these files¹¹ onto triples which we subsequently converted to

¹⁰https://www.mediawiki.org/wiki/Wikidata_Query_Service/User_Manual#Standalone_service

¹¹Described at https://amritasaha1812.github.io/CSQA/download_CQA/

COUNTLOGIC
Count Logical operators # Count Logical operators (Coreference)
UNIONMULTI
Union Multiple Relation
VERIFY3
3 entities, 2 direct, 2(direct) are query entities, subject is indirect # 3 entities, all direct, 2 are query entities

Table 8: Unseen question sub-types in SPICE non-i.i.d splits.

ttl format¹² with full URI to allow loading them to the KG query engine. We also filled missing information where possible, for example, missing relations such as “instance of” was filled with relation P31 and added data type information when this was omitted from the original files.

We used Blazegraph¹³ to deploy the local server, which uses only CPU-based resources and has access to 150G of RAM. Further details along with the script to host the server will be released upon acceptance.

C BertSP Model Configuration

Our model is implemented using pytorch (Paszke et al., 2019). For all experiments, we used the ADAM optimizer (Kingma and Ba, 2015) with 20,000 BERT warmup steps and 10,000 steps for decoder warm up. We use separate optimizers for the BERT encoder and decoder. BERT was fine-tuned during training with the initial learning rate set to 0.00002. A learning rate of 0.001 was set for the rest of model parameters. Our model was trained for 200,000 steps; we used 4 GPUs with 12GB of memory. We performed model selection on the validation set. We report results with the best performing model which had 6 decoder layers.

D Examples on Generalisation Splits

Table 9 shows examples from the generalisation splits: similar question sub-types seen during training, unseen question sub-type, and error on prediction. The most common error across different splits is that models use similar SPARQLs seen during training but fail to adapt them to the details (entities, types, relations, argument positions) in the unseen question sub-type. Other errors encompass using the incorrect SPARQL query (incorrect question intent) and incorrect entities and types.

¹²<http://www.w3.org/TR/turtle/>

¹³<https://blazegraph.com/>

COUNTLOGIC

SEEN	Union Single Relation	
	Which people are the creators of <i>The Theory of Everything</i> or <i>Ten Minutes to Live</i> ?	<pre>SELECT ?x WHERE { {wd:Q15079318 wdt:P162 ?x. ?x wdt:P31 wd:Q502895.} UNION {wd:Q7699260 wdt:P162 ?x. ?x wdt:P31 wd:Q502895.} }</pre>
UNSEEN	Count Single entity type	
	How many people starred in <i>Captain America: Civil War</i> ?	<pre>SELECT (COUNT(*) AS ?count) WHERE { wd:Q18407657 wdt:P161 ?x. ?x wdt:P31 wd:Q502895.}</pre>
UNSEEN	Count Logical operators	
	How many national association football teams or national sports teams represent Slovenia?	<pre>SELECT (COUNT (DISTINCT ?x) AS ?count) WHERE { {?x wdt:P1532 wd:Q215. ?x wdt:P31 wd:Q6979593.} UNION {?x wdt:P1532 wd:Q215. ?x wdt:P31 wd:Q1194951.} }</pre>
PRED		
		<pre>SELECT (COUNT(DISTINCT ?x) AS ?count) WHERE { {?x wdt:P1532 wd:Q215. ?x wdt:P31 wd: Q6979593 .} UNION {?x wdt:P1532 wd:Q215. ?x wdt: P31 wd: Q6979593.} }</pre>

UNIONMULTI

SEEN	Union Single Relation	
	Which people are the creators of <i>The Theory of Everything</i> or <i>Ten Minutes to Live</i> ?	<pre>SELECT ?x WHERE { ?x wdt:P915 wd:Q1247373. ?x wdt:P31 wd:Q838948.}</pre>
UNSEEN	Count Mult. entity type	
	How many people starred in <i>Django Kill</i> or <i>Shat-terday</i> ?	<pre>SELECT (COUNT(DISTINCT ?x) AS ?count) WHERE { {wd:Q1261875 wdt:P161 ?x. ?x wdt:P31 wd:Q502895.} UNION {wd:Q7490688 wdt:P161 ?x. ?x wdt:P31 wd:Q502895.} }</pre>
UNSEEN	Union Multiple Relation	
	Which administrative territories are the origin of <i>Les Chics Types</i> or are the native countries of <i>Robert Kuraś</i> ?	<pre>SELECT ?x WHERE { {wd:Q3231475 wdt:P495 ?x. ?x wdt:P31 wd:Q15617994.} UNION {wd:Q9310937 wdt:P27 ?x. ?x wdt:P31 wd:Q15617994.} }</pre>
PRED		
		<pre>SELECT ?x WHERE { {wd:Q3231475 wdt:P495 ?x. ?x wdt:P31 wd:Q15617994.} UNION {wd:Q9310937 wdt:P495 ?x. ?x wdt:P31 wd:Q15617994.} }</pre>

VERIFY3

SEEN	2 entities, both direct	
	Is <i>Zugspitze</i> located in Germany?	<pre>ASK {wd:Q3375 wdt:P17 wd:Q183.}</pre>
UNSEEN	3 entities, all direct, 2 are query entities	
	Is <i>Violet Oakley</i> a civilian of United States of America and Scheden?	<pre>ASK {wd:Q30 wdt:P27 wd:Q1226556. wd:Q557427 wdt:P27 wd:Q1226556.}</pre>
PRED		
		<pre>ASK {wd:Q1226556 wdt:P27 wd:Q30. wd:Q557427 wdt:P27 wd:Q557427.}</pre>

Table 9: Generalisation splits, unseen question sub-types, suport question sub-types seen during training, and example common erros on unseen predictions.

Table 10: Full list of question sub-types (intents) in SPICE. For each sub-type we show an example user question, whenever the question sub-type involves a conversational phenomenon (coreference or ellipsis); previous conversation interactions necessary for the interpretation of the user question are provided in grey.

Question Sub-Type	Example User Question
Clarification	
Simple Question Single Entity (Coreference)	U: Which political territory is that sporting event located in? S: Did you mean Speed skating at the 2010 Winter Olympics – Men's 500 metres? U: Yes
Logical Reasoning (All)	
Difference Multiple Relation	U: Which people were awarded with Order of Merit for Arts and Science and are not working as singer?
Difference Single Relation	U: Which international organizations had Poland but not Bulgaria as their member? U: Which city was Pierre Laffont born in? S: Marseille
Difference Single Relation (Ellipsis)	U: Which administrative territories are the sister cities of that city? S: Shanghai, Odessa, Naples U: But not Bologna
Intersection Multiple Relation	U: Which human settlements are situated close to Trave and have an adjacent border with Herzogtum Lauenburg?
Intersection Single Relation	U: Which works of art were filmed at Edinburgh and Berlin? U: Which language does José María Lassalle speak in? S: Spanish
Intersection Single Relation (Ellipsis)	U: And also Sergio Gil
Union Multiple Relation	U: Which watercourses are located in the neighbourhood of Bremen or are the tributaries of Ob?
Union Single Relation	U: Which people are the creators of The Theory of Everything or Ten Minutes to Live?
Union Single Relation (Ellipsis)	U: What is the profession of Mai Yamada? S: announcer U: Or Kazimierz Rogoyski?
Quantitative Reasoning (All)	
Min/Max Single entity type	U: Which musical instruments are played by min number of people?
Min/Max Mult. entity type	U: Which organizations are the main building contractors of max number of architectural structures and buildings?
Atleast/ Atmost/ Approx. the same/Equal Single entity type	U: Which musical instruments are played by exactly 5 people?
Atleast/ Atmost/ Approx. the same/Equal Mult. entity type	U: Which events are demonstrated in atleast 3 prints and genres of sculpture?
Comparative Reasoning (All)	
More/Less Mult. entity type	U: Which landforms are known for containing lesser number of chemical compounds or minerals naturally than Stetind pegmatite?
More/Less Mult. entity type (Ellipsis)	U: Which landforms are known for containing lesser number of chemical compounds or minerals naturally than Stetind pegmatite? S: Euboea, Izalco, Mount Nyiragongo U: And also tell me about Tufte quarry?
More/Less Mult. entity type (Coreference)	U: Which administrative territory is that person a civilian of? S: Spain U: Which administrative territories are the countries of origin of lesser number of television programs or works of art than that administrative territory?
More/Less Single entity type	U: Which television programs have been dubbed by more number of people than Puss in Boots: The Three Diablos?
More/Less Single entity type (Ellipsis)	U: Which television programs have been dubbed by more number of people than Puss in Boots: The Three Diablos? S: House, K-On!, K-On!! U: And also tell me about Chip 'n Dale Rescue Rangers?
More/Less Single entity type (Coreference)	U: Which languages are max number of literary works composed in? S: English U: Which languages are the mother tongues of less number of people than that language?
Simple Question (Direct)	
Simple Question	U: Which type of sport did Amel Tuka participate in?
Single Entity	U: What is the capital of Sweden?
Mult. Entity	U: Who were the writers of On being and essence, De vegetabilis et plantis libri septem and Historia de regibus Gothorum, Vandalorum et Suevorum?
Simple Question (Ellipsis)	
only subject is changed, parent and relation remains same	U: Which organizations are the sponsors of Janice Anderson? S: Montrail, Patagonia, Inc. U: And also tell me about Manikala Rai?
object parent is changed, subject and relation remain same	U: Which watercourses are situated nearby Munich? S: Eisbach, Würm, Isar U: And which river?
Simple Question (Coreference)	
Mult. Entity	U: Which releases have Motown as their record label? S: What's Going On, Got to Be There, Can't Slow Down U: Which genre do those releases belong to?
Single Entity (Coreference)	U: Which narrative location is The Penalty set in? S: San Francisco U: Which color is associated with that film genre?
Verification (Boolean) (All)	
2 entities, both direct	U: Is Zugspitze located in Germany?
2 entities, one direct and one corefered, object is corefered	U: Which university was Eden Stiles educated at? S: University of Michigan U: And what about C. V. Raman?

Continued on next page

Table 10 – Continued from previous page

Question Type/Sub-Type	Example User Question
	S: University of Madras U: Was Ravindra Wijegunaratne educated at that university?
2 entities, one direct and one corefered, subject is corefered	U: Which German business organization was Gustav Peichl a member of? S: Academy of Arts, Berlin U: What was designed by that person? S: Millennium Tower U: Does that tower block belong to Austria?
3 entities, 2 direct, 2(direct) are query entities, subject is corefered	U: Which administrative territory was Gary Collier born at? S: Fort Worth U: Is that administrative territory a sister city of Adamsville, New Brunswick and Yuen Long Kau Hui?
3 entities, all direct, 2 are query entities	U: Is Aix-en-Provence partner town of Baton Rouge and Hemmatabad, Alborz?
one entity, multiple entities (as object) coreferred	U: Which armed conflicts are Battle of the Argeş or Battle of the Yellow Sea a part of? S: Romania during World War I, Russo-Japanese War U: Did those armed conflicts fight in Rui Natsukawa?
Quantitative Reasoning (Count) (All) Incomplete count-based ques	U: How many people influenced Chris Marker? S: 1 U: And also tell me about Ada Yonath? S: 1 U: And what about Mikhail Bakunin?
Count over Atleast/ Atmos/ Approx. the same/EqualMult. entity type	U: How many cities are the terminus locations of atleast 5 thoroughfares and roads?
Count over Atleast/ Atmos/ Approx. the same/EqualSingle entity type	U: How many musical instruments are played by exactly 2 people?
Count Logical operators	U: How many bodies of water or watercourses are situated nearby Lübeck?
Count Logical operators (Coreference)	U: Which administrative territory is the native country of Carolina Goic Boroevic? S: Chile U: Who is the head of the government of that administrative territory? S: Michelle Bachelet U: What is the capital of that administrative territory? S: Santiago U: How many capitals or cities are sister towns of that city?
Count Mult. entity type	U: How many people starred in Django Kill or Shatterday?
Count Single entity type	U: How many people starred in Captain America: Civil War?
Count Single entity type (Coreference)	U: Which armed conflict did Lionel of Antwerp, 1st Duke of Clarence take part in? S: Hundred Years' War U: How many people did that armed conflict engage in?
Comparative Reasoning (Count) (All) Count over More/Less Mult. entity type	U: How many administrative territories have adopted lesser number of holidays and people as patron saint than Santo Stefano al Mare?
Count over More/Less Mult. entity type (Ellipsis)	U: How many administrative territories have adopted lesser number of holidays and people as patron saint than Santo Stefano al Mare? S: 296 U: And what about San Donato Milanese?
Count over More/Less Mult. entity type (Coreference)	U: Which administrative territories are Luigi Einaudi the head of state of and have UTC+01:00 as their time zone? S: Italy U: How many administrative territories are the origins of greater number of literary works or releases than that administrative territory?
Count over More/Less Single entity type	U: How many legislatures represent lesser number of states than East Bengal Legislative Assembly?
Count over More/Less Single entity type (Ellipsis)	U: How many legislatures represent lesser number of states than East Bengal Legislative Assembly? U: 207 U: And how about Estates of Curaçao?
Count over More/Less Single entity type (Coreference)	U: Which french administrative division was Philippe Esnault born in? S: Alençon U: Which occupation has that person as his/her's career? S: historian U: Which administrative territory is the native country of that person? S: France U: How many administrative territories inspired less number of fictional locations than that administrative territory?